



Data Release Use Case Team: Status of effort and issues identified

September 2013

Team Leads: Fran Lightsom and Viv Hutchison

Presenter: Keith Kirk

U.S. Department of the Interior
U.S. Geological Survey

Serving USGS
Information to the Nation

The mission of the USGS
is to conduct science and to
make available the resulting
knowledge and information.
USGS has a long history of
service to the Nation.

and virtually any...
that will identify any...
make our... sure a...
knowledge at any...
where at any...
this desired...
state, the USGS will...
to serve as the Nation's...
natural science...

Members of the CDI Use Case Team

- Fran Lightsom (co-lead), Natural Hazards (Woods Hole, MA)
- Viv Hutchison (co-lead), Core Science Systems (Denver, CO)
- John Faundeen, Climate and Land Use Change (Sioux Falls, SD)
- Greg Gunther, Energy and Minerals (Denver, CO)
- Keith Kirk, Office of Science Quality and Integrity (Santa Cruz, CA)
- Greg Miller, Natural Hazards (St. Petersburg, FL)
- Andrea Ostroff, Core Science Systems (Reston, VA)
- Carolyn Reid, Office of Science Quality and Integrity (Reston, VA)

- Facilitator: Peter Fox, Rensselaer Polytechnic Institute (Troy, NY)

*****Cross-Mission Area Representation*****

Background: Use Case Team

- **Original purpose:** to develop a process, based on current policies and workflows, that enables USGS employees to determine if a particular set of data is approved for release
- Convened early 2012 through the USGS Community for Data Integration (CDI) Data Management Working Group
 - Met face-to-face in Reston VA in April 2012
 - Weekly phone meetings ever since (*whew!*)



Major Challenges Initially Identified:

- **Lack of:**

- bureau-wide understanding of policies and procedures for releasing data
- bureau-wide understanding of distinctions between publishing data in a USGS series report versus other means of data release
- attention to data preservation within Fundamental Science Practices
- Explanation of differences between peer review and data review is not reflected in current policy

- **Resistance to:**

- metadata creation along with inconsistent or absent treatment of metadata in the release process



Connections to External Drivers

- **Open Government Initiatives supporting broader public access to Federal and Federally-supported data and information**
 - *the Use Case Team thinking was ahead of recent directives, now we have the opportunity to leverage these directives to facilitate positive change in the USGS*

Policy	Type	Date Issued
Transparency and Open Government	Presidential Memorandum	1/21/2009
Open Government Directive	OMB Memorandum M-10-06	12/8/2009
Digital Government: Building a 21st Century Platform to Better Serve the American People	Federal CIO Strategy Document	5/23/2012
Managing Government Records Directive	OMB-NARA Memorandum M-12-18	8/24/2012
Increasing Access to the Results of Federally Funded Scientific Research	OSTP Memorandum	2/22/2013
Making Open and Machine Readable the New Default for Government Information	Executive Order 13642	5/9/2013
Open Data Policy-Managing Information as an Asset	OMB Memorandum M-13-13	5/9/2013

Web Release Use Case: Diagrams

Data Release via Web: Assumptions

To build the Use Case, some assumptions were made to start:

- Data is not interpretive
- USGS data product is assumed to be non-proprietary, and non-sensitive
- Science Center web sites can host data

Data Release via Web: Pre-Conditions

To build the Use Case, we noted existing “pre-conditions” before someone would use this workflow:



- USGS data exists and is available to the Author
- Data has not been previously released nor is part of a national collection or other ‘approved’ dissemination
- Data are not pre-decisional (SM 502.5)
- A web site that is appropriate for release of the data exists or can be created and is available to the Author.
- Data and datasets are non-interpretive and therefore require a *data* quality review not a traditional peer review.
 - Peer review is appropriate for interpretive information products (scholarly publications).

Actors

- A use case establishes the relevant “actors”...what roles are involved?



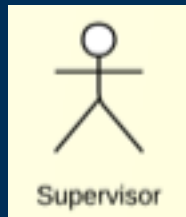
creates and revises product, initiates product approval request, prepares metadata for product



reviews product for scientific quality



manages product preservation, generally this is the originating Science Center or NatWeb



reviews for conformity with FSP policies and processes



approves product for release



oversees application of digital object identifier for data and metadata



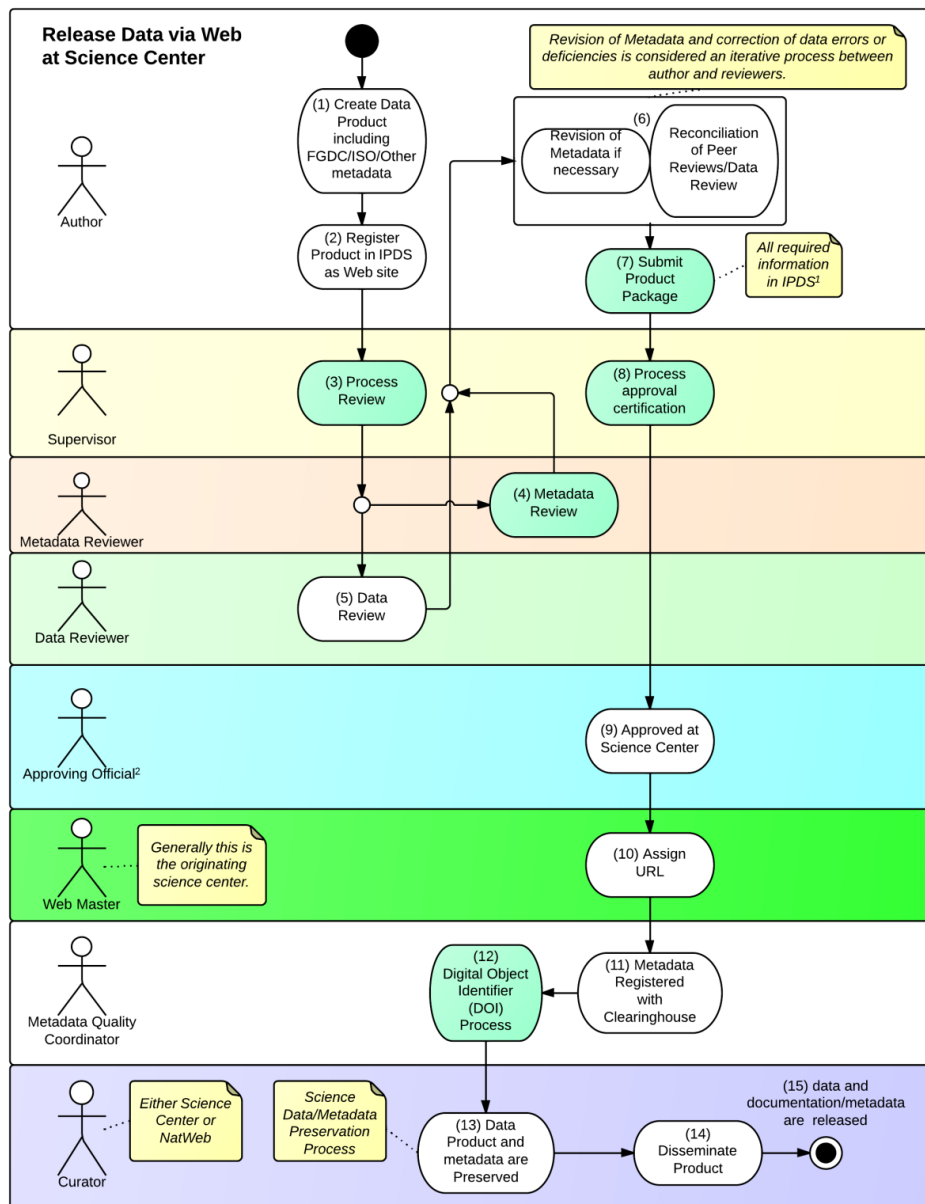
reviews metadata for accuracy and conformance with standards



generally this is the originating Science Center

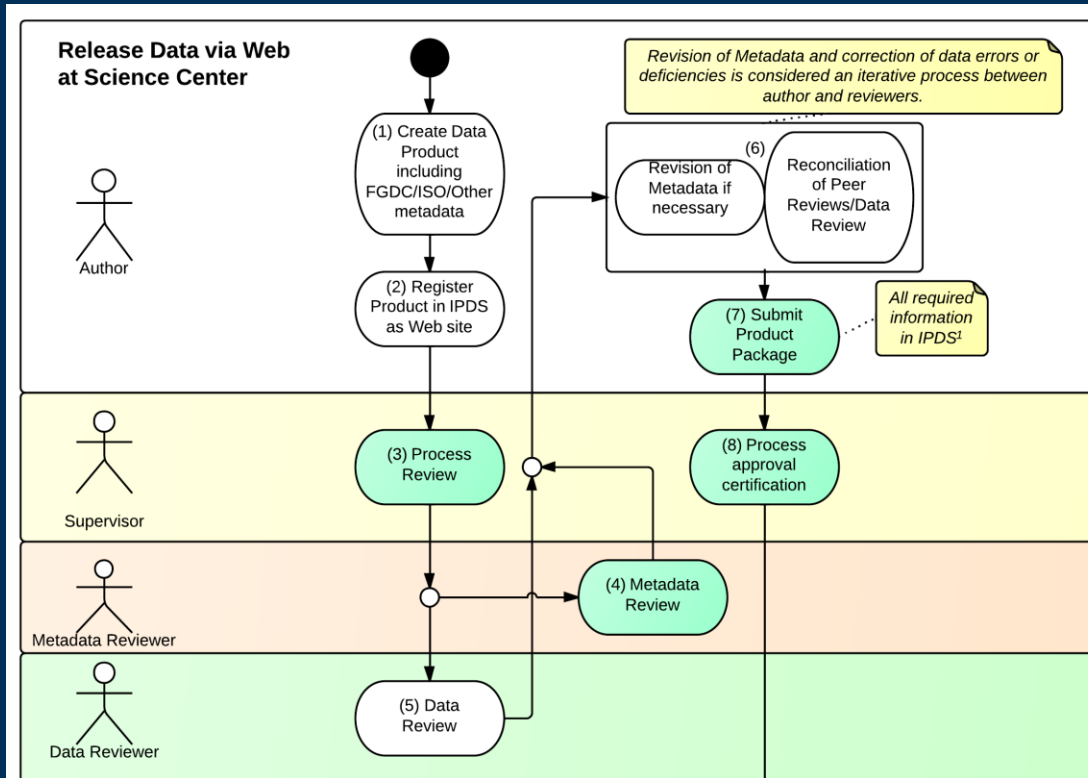
Data Release via Web:

Workflow Diagram Overview



1. IPDS is the Bureau Information Product Data System. Refer to URL: <http://internal.usgs.gov/publishing/ipds.html>
2. Under FSP approval of non-interpretive information products is delegated to the Science Center Director

Use Case: The details

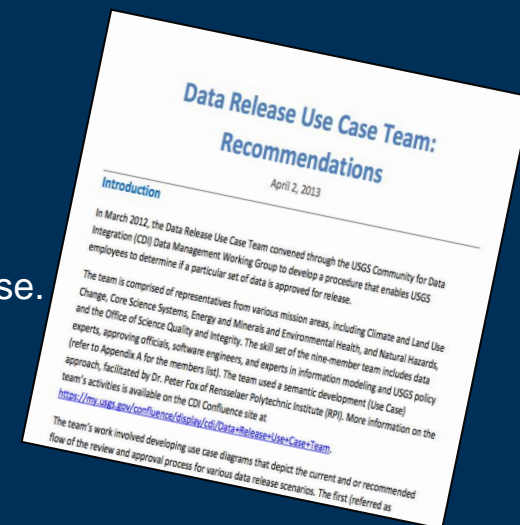


[Link to PDF:](#)

What did we learn?

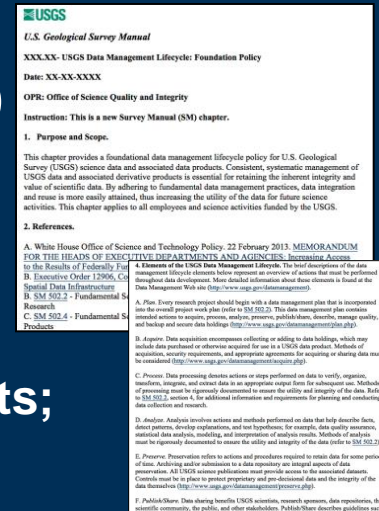
Use Case Recommendations

1. **Implement Digital Object Identifiers** or DOIs, to enable identification of both data that are approved for release and the associated publications.
2. Search capabilities for **publication catalogs that link to data**. (Possible with new Pubs Warehouse, ScienceBase2)
3. Implement USGS online data services that **provide data in a useful format**.
4. **Training** for providing data in useful formats on web.
5. **Guidance** for choice of data reviewers.
6. **Guidance** on responsibilities related to data review.
7. Hold data until data is reviewed and metadata is appropriate.
8. Require a **DOI in metadata record** to identify data that is approved for release.
9. **Enforce metadata review** for data that is released.
10. Education and resources for metadata reviewers.
11. Guidance for determining interpretive content.
12. **Policies and guidance for data preservation** (FSPAC working on this).
13. Develop approved **online repositories** that will preserve USGS data and information and ensure that they can be found and used in the future.
14. **Mandatory** training on data release (i.e., DOI Learn).
15. Provide information about new policies and processes for web release of data.
16. Establish **RGE credit** for high-quality data release.
17. Establish **data citation standards** so scientists get credit.



Connecting the Dots in USGS policy

- In Progress
(via the “Data Mgmt Policy Group” & FSPAC)
 - Data Management Overarching Policy
 - Based on USGS Data Lifecycle
 - Metadata for Datasets and Information Products
 - Requires FGDC endorsed standards for data products; use of IPDS metadata for publications
 - Release of Software
 - Release of Data
 - Requires application of Digital Object Identifiers (DOIs) for all data (to connect to publications)
- Other policies needed:
 - Data citation standards
 - RGE recognition so that scientists get credit for releasing data
 - Data preservation standards



Build It, Test It

- **Washington Water Science Center**

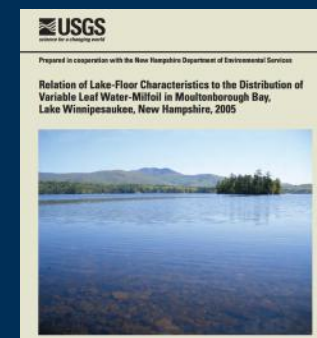
(First data set with DOIs released 6/10/2013 in response to litigation associated with Elwha Dam Removal). Data released initially to meet immediate need to provide data in support of litigation against Interior related to consequences of Elwha River dam removal effort. Ongoing data collected daily is added to site after QA/QC

- **Pacific Science Center: Santa Cruz**

Rich and VERY large data set. Not appropriate for Data Series owing to data type. Large data stacks using Open Access. YouTube and Picasso leveraged as data repositories. Separate YouTube channel approved by OCAP for this data. Approved using new IPDS. Currently served on non-USGS site will be moved to USGS site as resources allow. This will be easily facilitated by DOI metadata. Metadata for all video and images per OCAP requirements.

- **Woods Hole**

Data were used in figures in a 2007 Scientific Investigations Report (SIR). The data were not published at the time but now requests for data are numerous. Data will be released outside of USGS series to facilitate need for machine readable data. Data will go online with a metadata record after an FSP-compliant review and approval process. Data will probably be part of the existing Woods Hole Science Center data library Website.



Benefits of Web Release for USGS Data

- **Eliminates costs associated with producing and distributing media**
 - uses existing IT resources
- **Allows science center directors to use discretion in designating materials for editorial review**
 - more efficient use of the editors' time and a higher rate of information product publication.
- **Open Data Compliance: Data online can be hosted in formats that allow re-use**
 - Moving away from a pdf format to open data formats will greatly enhance our ability to share data and make our data more interoperable
 - Compliance with OSTP memo and other directives.
- **Helps meet OSTP & OMB Directives related to open access, data release & Machine Readability**



Next Steps for Data Release Team

- **Seek review and approval of Use Cases:**
 - by broader groups (CDI DMWG, SDCN, etc)
 - by Alan and Kevin and USGS ELT
- **Post Use Cases on USGS Data Management website**
- **Develop:**
 - more use cases
 - more recommendations
 - education components
 - **Ensure relevant policies and processes are developed (i.e., data release policy, digital object identifier guidance, etc)**
- **Get the word out – make processes clearer and easier**



**KEEP
CALM
AND
RELEASE
DATA**

Thank you!

Questions and Comments



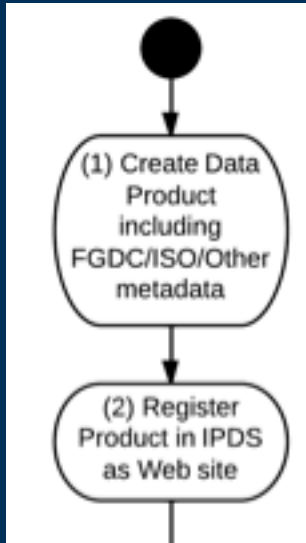
**Oh! Did you say you want to review
use case materials? 😊**

- <https://my.usgs.gov/confluence/display/cdi/Web+Release+Use+Case>
- **send comments to Fran (flightson@usgs.gov) and Viv (vhutchison@usgs.gov)**

Process Details



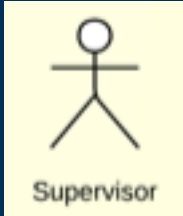
creates and revises product, initiates product approval request, prepares metadata for product



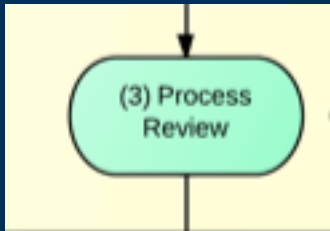
Author develops Product:

- Identifies the appropriate Web site (or Web data service)
- Makes available for review:
 - one or more clean data sets in appropriate format;
 - one or more metadata records;
 - any additional descriptive materials needed to ensure the data are discoverable and useful;
 - draft Web site

Process Details



reviews for
conformity with
FSP policies
and processes



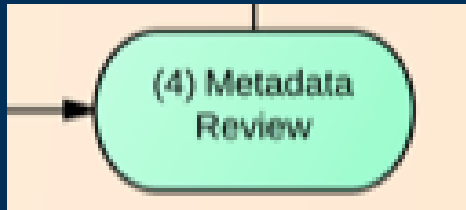
Process Review:

- Supervisor verifies appropriateness of choice of Web release;
- ensures the draft Web site meets standards for completeness,
- verifies choice of data reviewers, metadata reviewers;
- ensures compliance with relevant policy requirements
- gives author permission to provide draft web site to reviewers.

Process Details



reviews metadata for accuracy and conformance with standards



Metadata Review:

- Produces written report and returns this report to the author.

• Metadata Review Process:

- Check compliance using a recommended metadata validation tool.
- Perform metadata quality checks:
 - Check that the metadata matches the data
 - Check that data field names and values are defined and consistent with information in entity/attribute section of metadata record
 - Check that bounding coordinates match location keywords
 - Check temporary on-line linkage to data exists (this link(s) will change when final DOI is assigned)
 - Check that information about processing steps, methodology, lineage are included in the record.
 - Does the metadata provide robust information about how to use the data files – access instructions, software requirements, data models, definitions of terms, and so on?

Process Details



reviews product for
scientific quality

Data Review:

- Produces written report to the author with a recommendation that Product be released and a list of any recommended changes.

Data review may include the following:

- Is data format reasonable for public distribution (released using common standards)?
- Are data values reasonable? Are they in a valid range for that measurement, do they display any expected seasonal or daily trends, is there consistency between adjacent or otherwise related datasets?
- Can data be used by appropriate analysis and visualization tools?
- Does the metadata match the data?
- Are the data attributes listed in the metadata in agreement with the data?
- Are the techniques and methods scientifically sound and well described? Could a knowledgeable scientist or technician recreate the final data set from the descriptions? Can this information be easily found and used?
- Does the Product as a whole, through its design or documentation, provide enough information that the data and metadata can be easily found and used

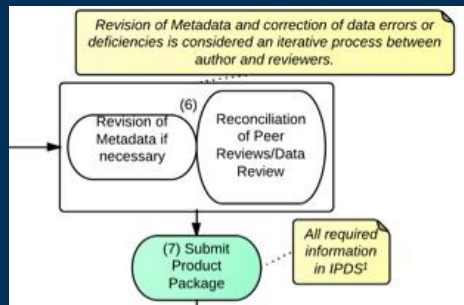


Process Details



creates and revises product, initiates product approval request, prepares metadata for product

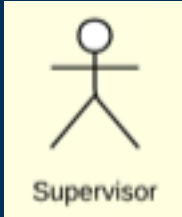
- Revision Process:
 - Author iteratively revises Product in response to and in cooperation with metadata reviewer and data reviewer comments and documents reconciliation and responses.
- Submit Package:
 - Author places all relevant materials in IPDS document vault as verification the review/reconciliation of the product took place and notifies Supervisor.



These materials include:

- A link to the draft Web site which has been reviewed and revised.
- Names of data reviewers and metadata reviewer, and evidence that they agree with revised form of the data and metadata.
- Reports from data and metadata reviewers, annotated by author to indicate changes made to data and metadata in response to reviews.
- Information about versioning, if the data is a first version that must be clearly indicated.
- URLs, DOIs, or bibliographic citations of publications that are related to the data.

Process Details



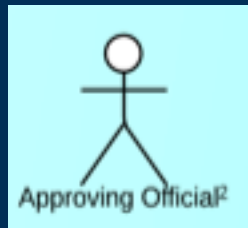
reviews for conformity with
FSP policies and
processes



Process Approval Certification:

- Supervisor verifies that the draft Web site is complete and responses/reconciliations are appropriate. Supervisor forwards request to local Approving Official (Science Center Director).

Process Details



approves
product for
release



Approval:

Non-interpretive data product is approved at science center or rejected.

(Science Center Director = Approving Official in case of non-interpretive data)

If rejected author is informed of need for additional changes.

Process Details



generally this is
the originating
Science Center

URL is assigned:

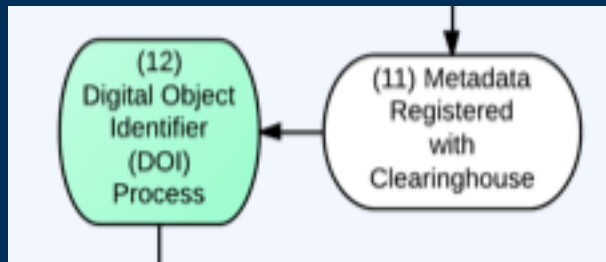
- The URL assigned at this step represents the physical/actual URL
 - (EXAMPLE:
<http://energy.usgs.gov/data/dataset.zip>)
 - URL is used as the de-referencable URL when obtaining the DOI



Process Details



oversees
application of digital
object identifier for
data and metadata

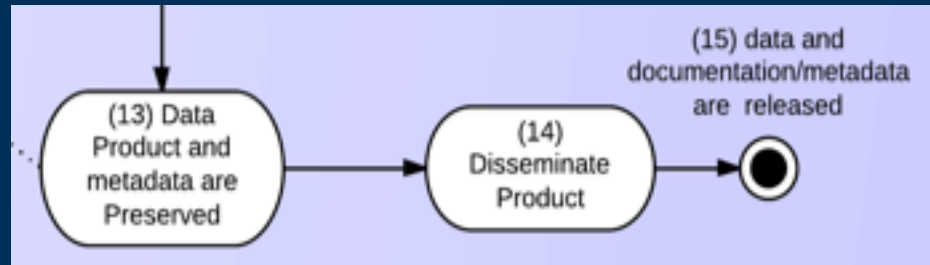


- **Digital Object Identifier obtained for data and populated into all relevant metadata, FGDC/ISO**
 - **Note: Updated metadata with new DOI (at least the DOI URL) must be passed back to author or responsible contact at Science Center**
- **Metadata Registered with Core Science Metadata Clearinghouse to make it available, according to Executive Order 12906 (1994)**

Process Details....final stretch



manages product
preservation, generally
this is the originating
Science Center or
NatWeb



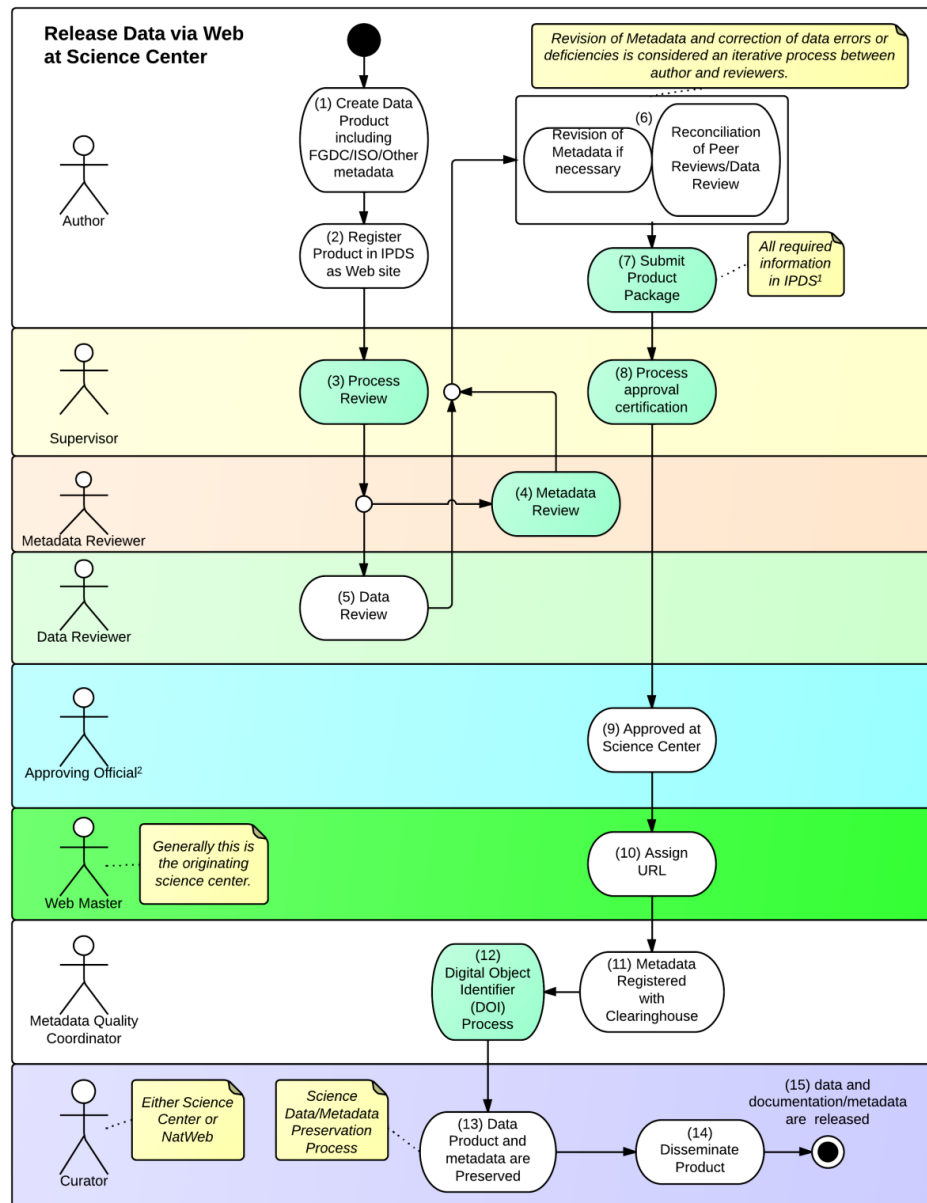
Preservation:

- **Data Product and metadata are preserved**
 - If product is on NatWeb, then archive requirements are met
 - If product is being made available elsewhere, ensure both the metadata and the data are preserved according to standards.

Data Release!

- **Approved Product is disseminated by Web release at Science Center.**

Data Release via Web

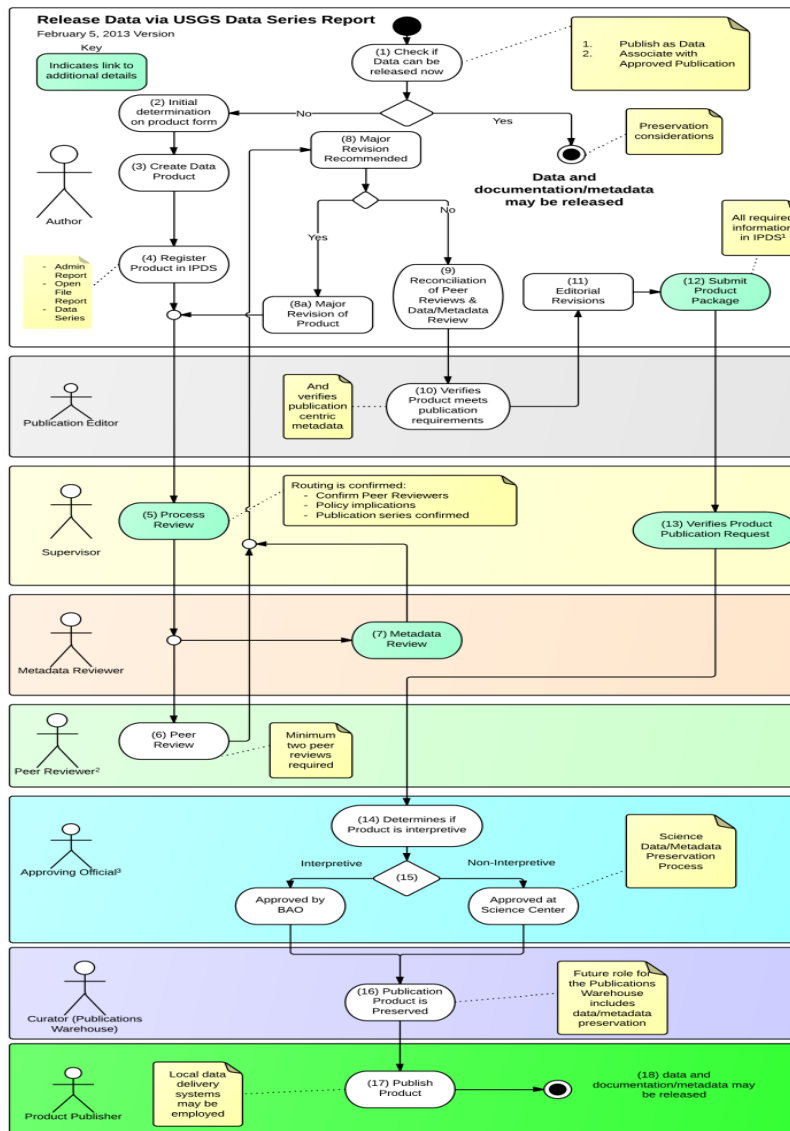


1. IPDS is the Bureau Information Product Data System. Refer to URL: <http://internal.usgs.gov/publishing/ipds.html>
2. Under FSP approval of non-interpretive information products is delegated to the Science Center Director

Additional Slides...

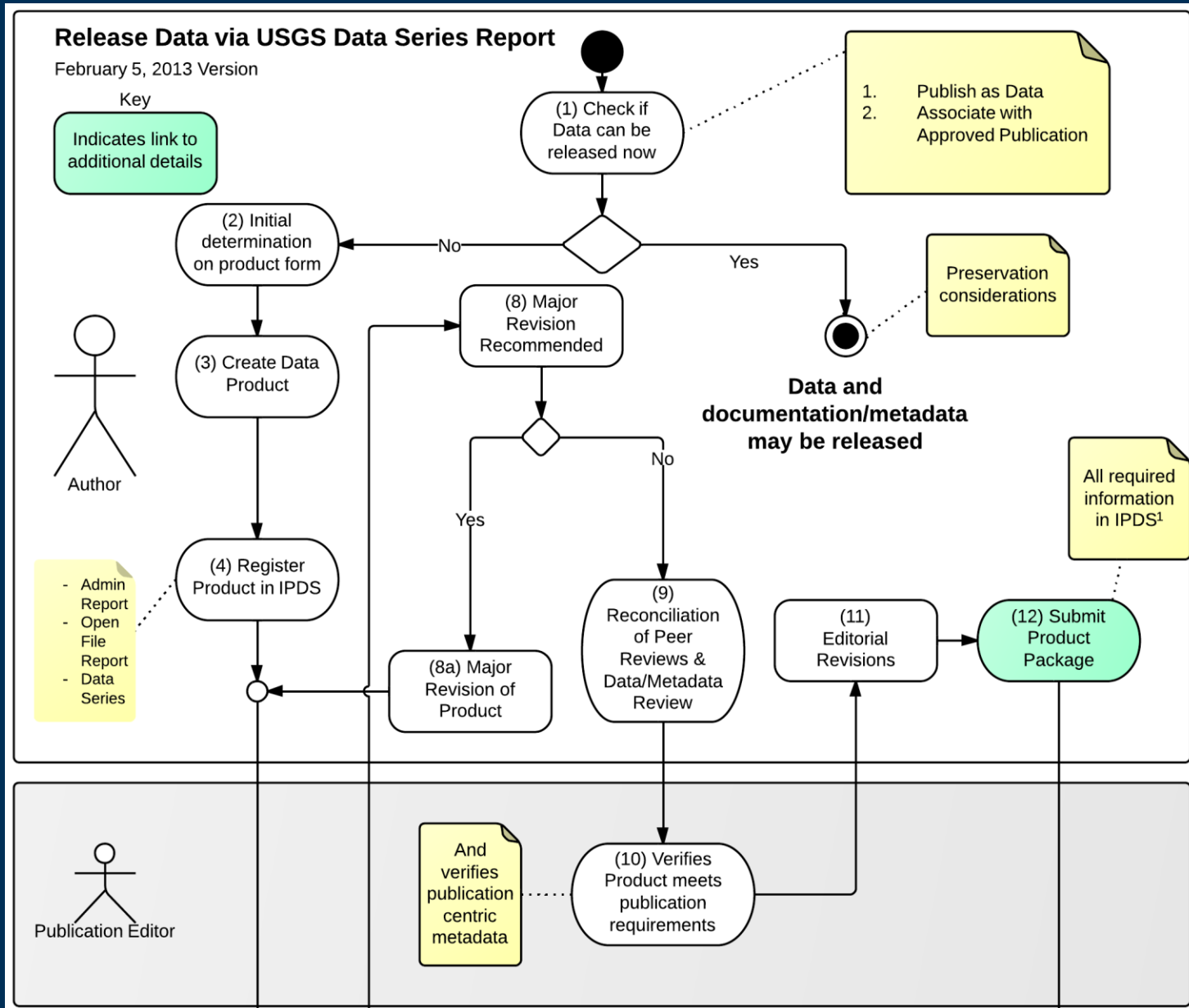
Data Release through USGS Series Publication

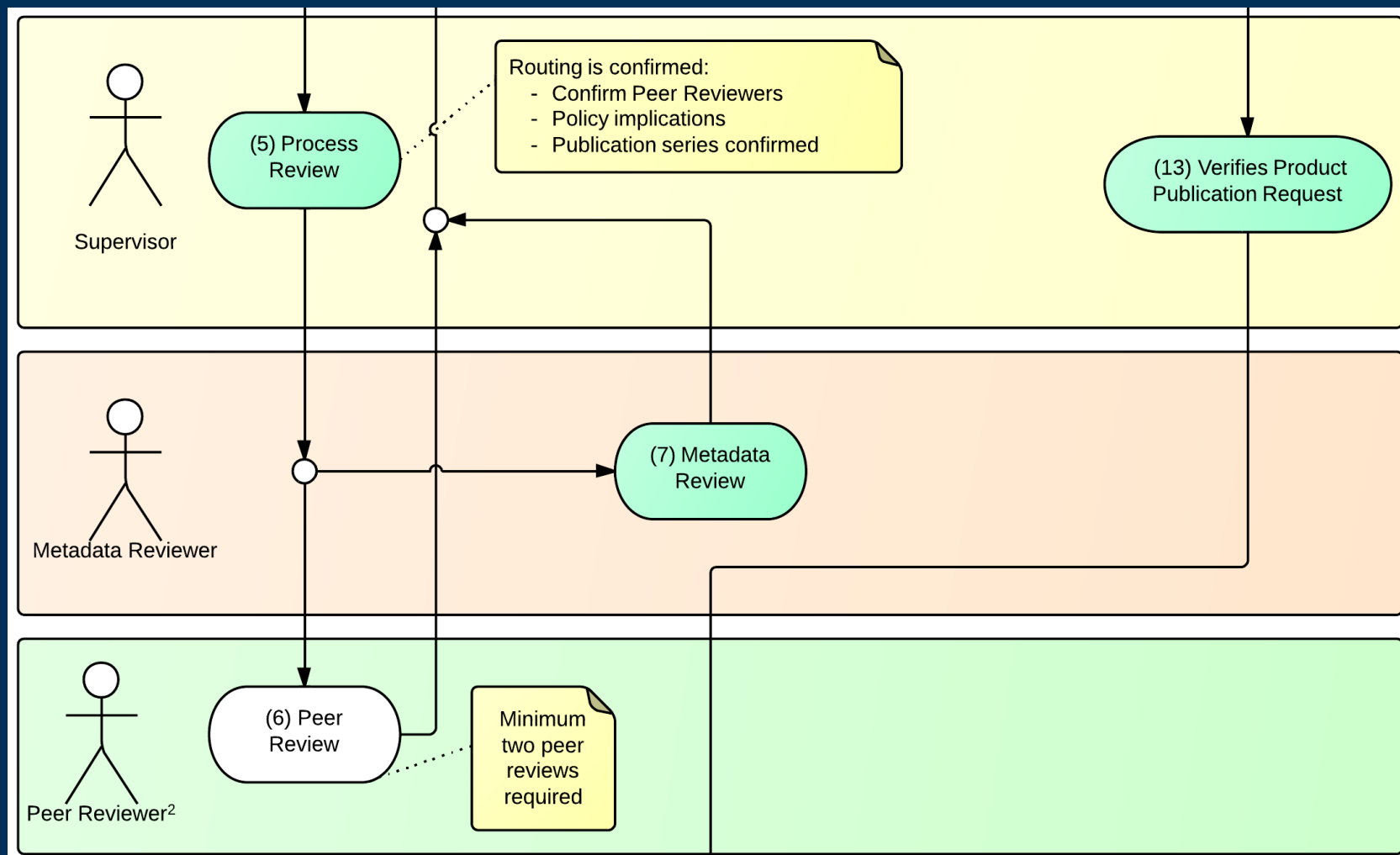
DATA SERIES USE CASE DIAGRAMS

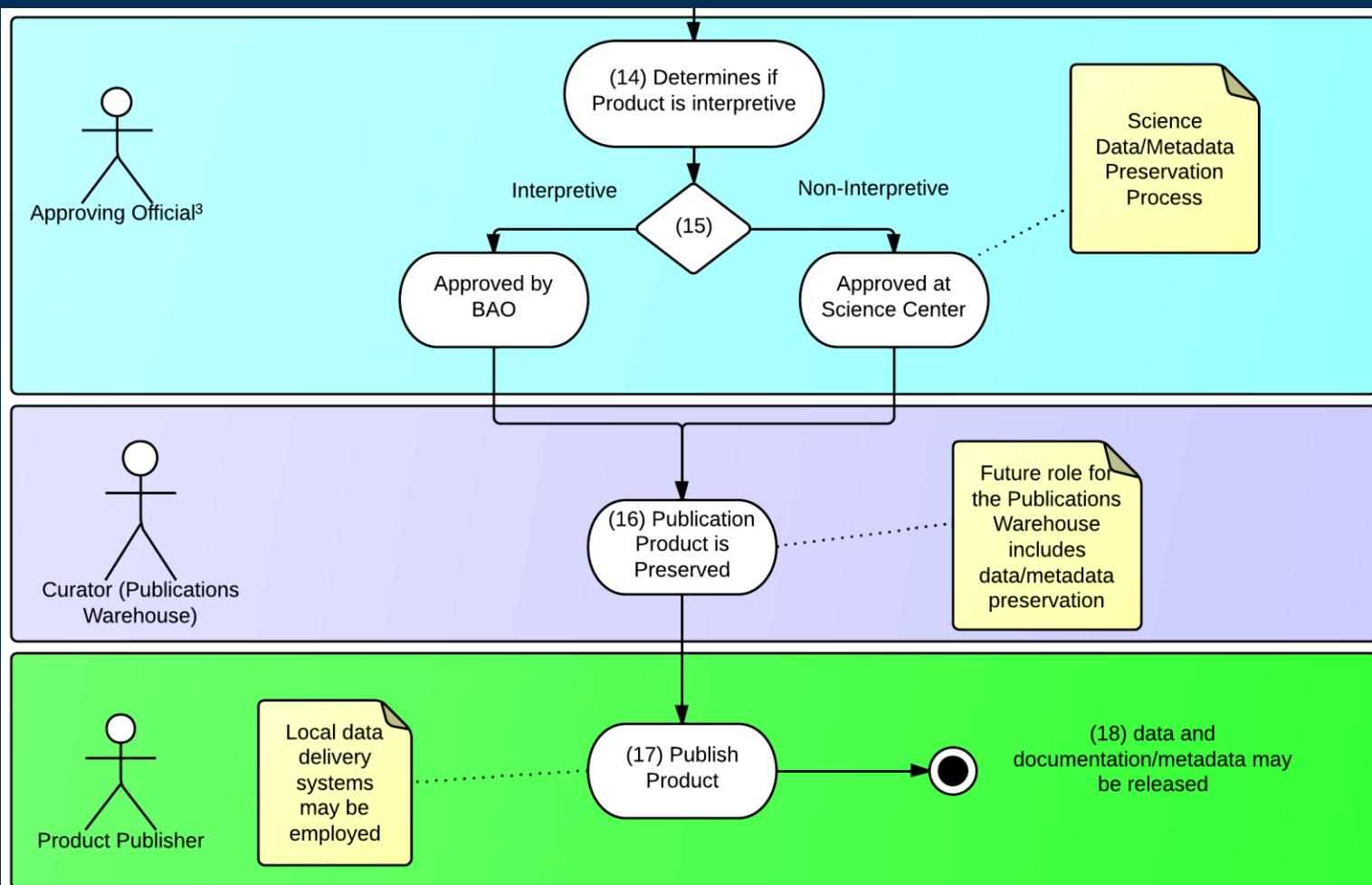


Release Data via USGS Data Series Report

February 5, 2013 Version







1. IPDS is the Bureau Information Product Data System. Refer to URL: <http://internal.usgs.gov/publishing/ipds.html>
2. If reviewer(s) recommend that Product not be published, Author performs major revision to Product, possibly including extensive clean-up of data and/or metadata, and starts again at flow step 5
3. The Science Center Director or designee determines if product is interpretive or non-interpretive. If non-interpretive may approve. In addition, either the Science Center Chief or the BAO and "Kill" the process requiring additional information or a restart. Author performs major revision to Product, possibly including extensive clean-up of data and/or metadata, and starts again at basic flow step 2